# A Framework for Scalable Trainable Image-based Query in Video

Jianbo Shi

Informedia-II
*Carnegie Mellon University*
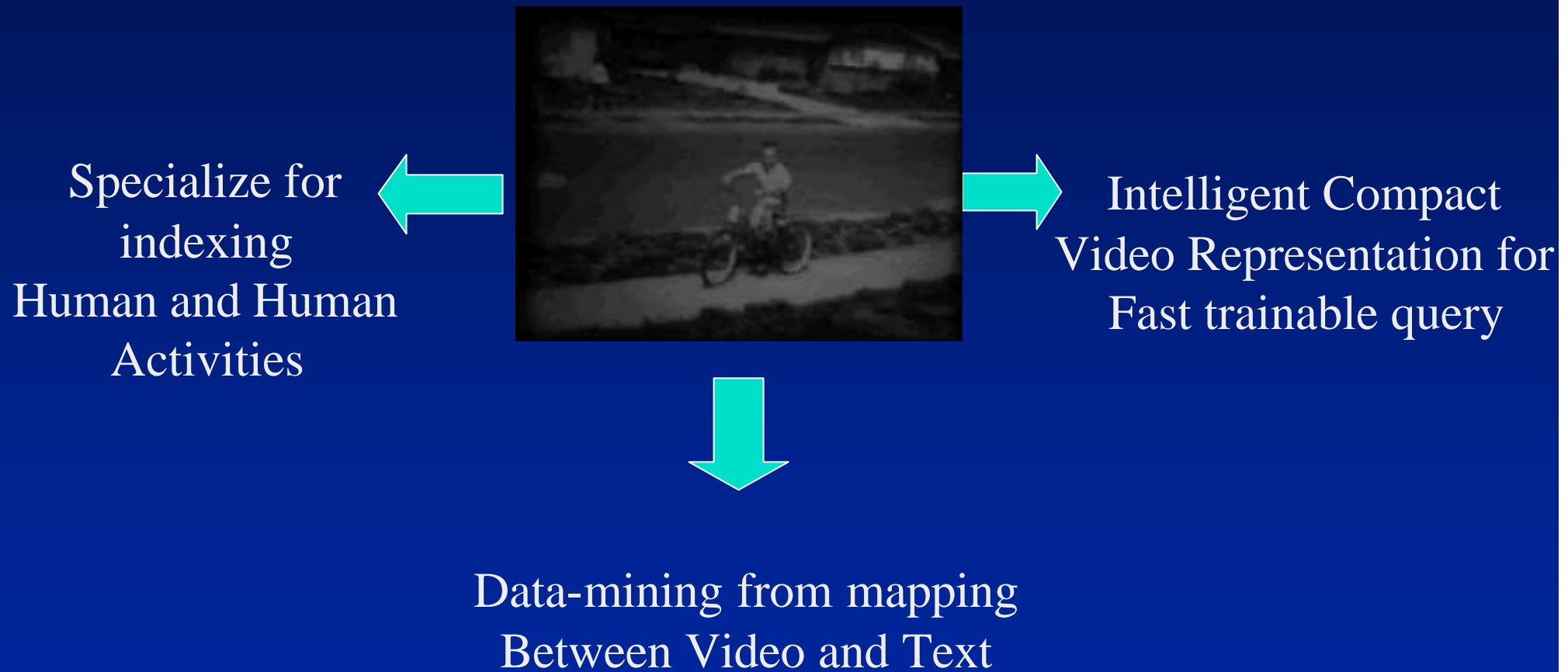
*informedia*

- **Multi-modal query of video**
  - Speech recognition
  - Text analysis
  - Object detection in video

- **Visualization and Summarization**
  - Multiple Video Documents
  - Topic collages

http://www.informedia.cs.cmu.edu

- Image-based query is a powerful tool for finding relevant information
  - Informedia-I, Blobworld, QBIC, ect.

*informedia*

- **Image-based query is a powerful tool for finding relevant information**
  - Informedia-I, Blobworld, QBIC, ect.

- **…but we are far from achieving**
  - Human level object recognition
  - Rapid processing of massive video/image data
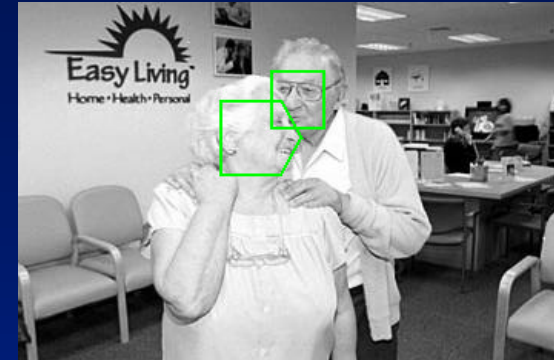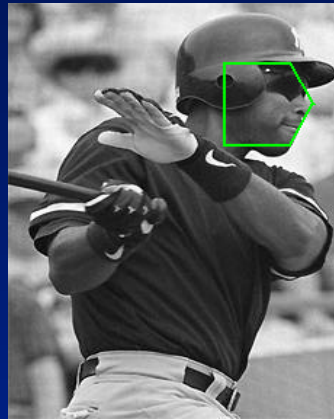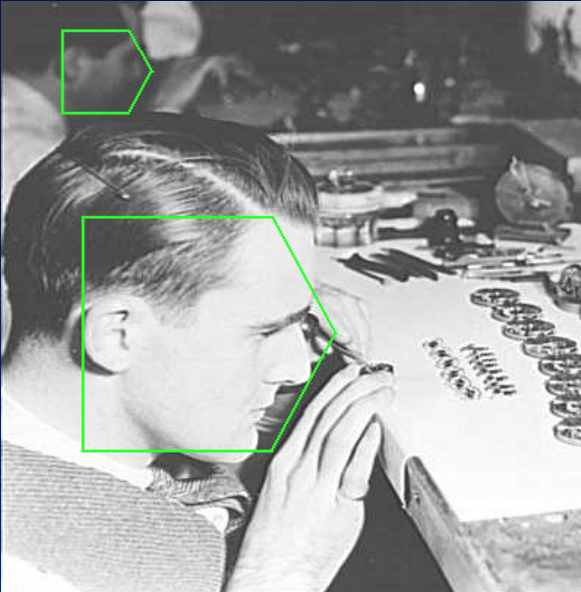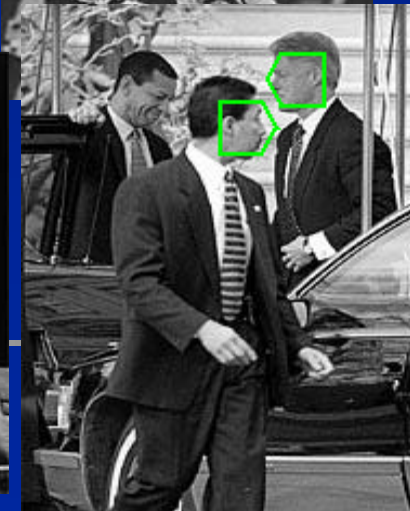    - *The story is in the large of volume data*

**Specialize for indexing Human and Human Activities**

**Intelligent Compact Video Representation for Fast trainable query**

**Data-mining from mapping Between Video and Text**
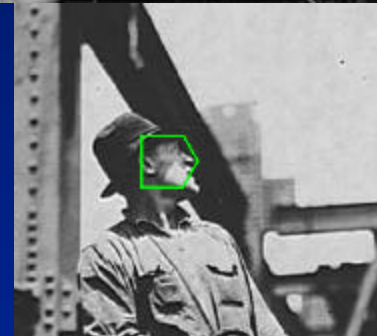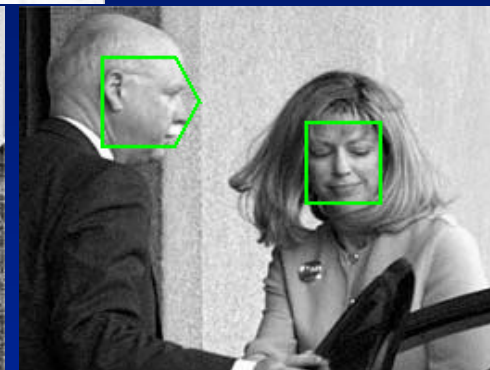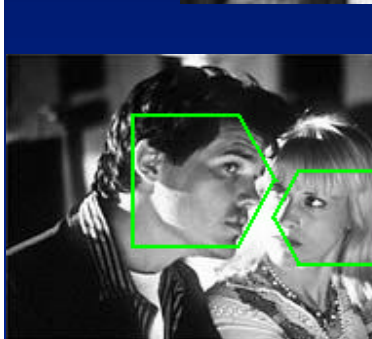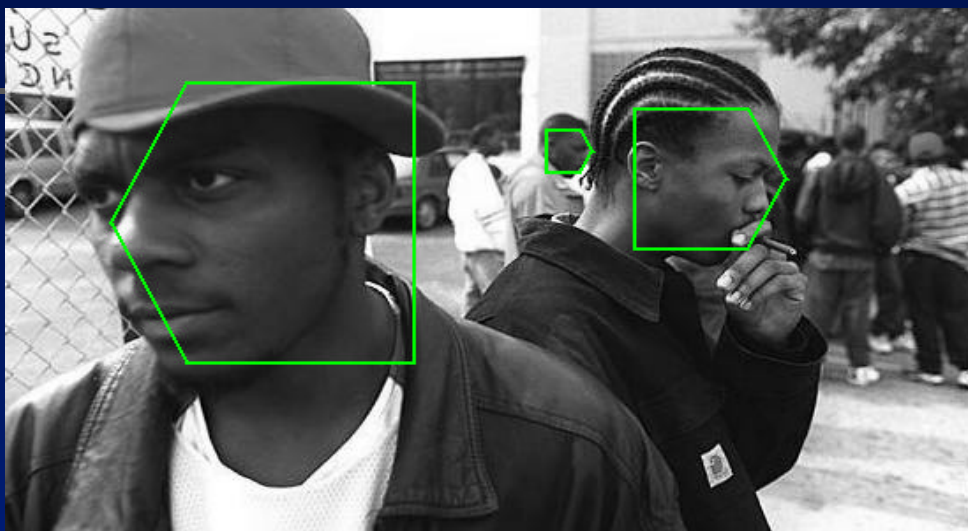
Human Identity: Face



Human Activity Recognition: Body

Face Detection - Henry Schneiderman (CMU)

Detecting detecting human in video

Understanding human activity

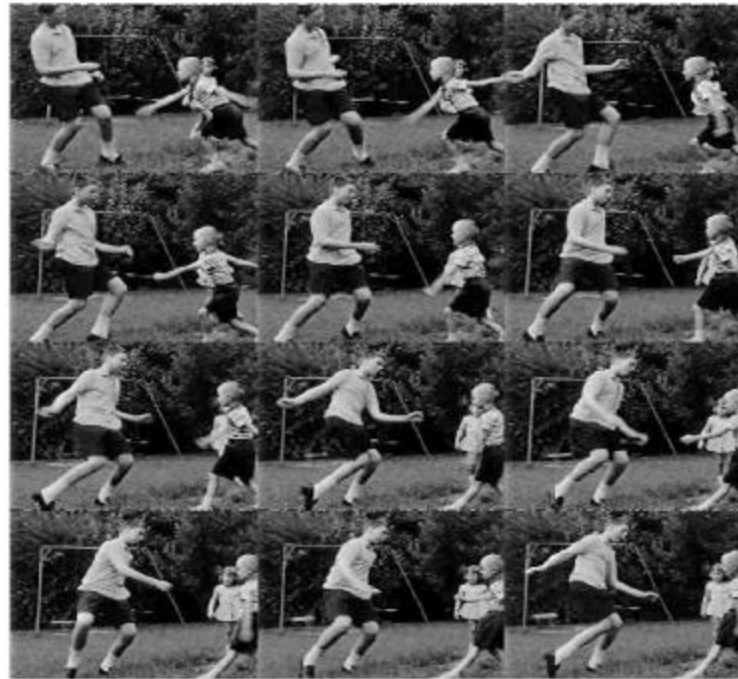Measuring/analyzing body movement

Recognizing human from body movement
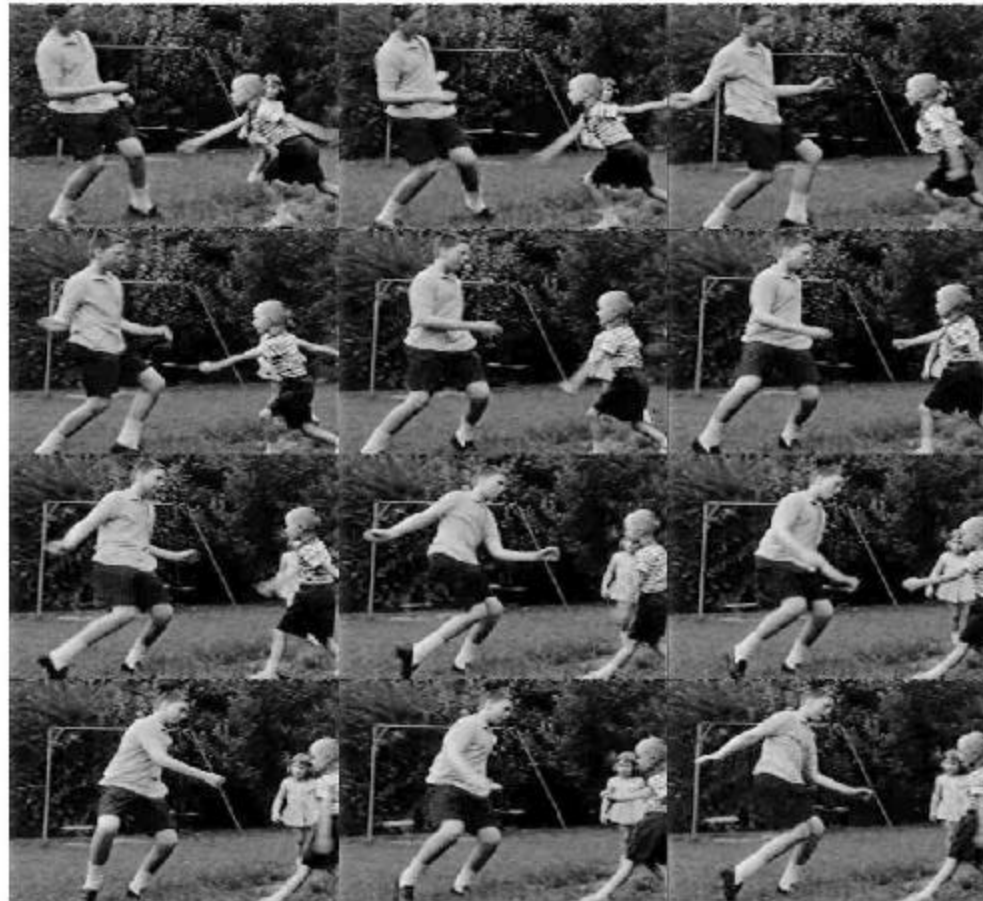
$$N = (npixel \times num\_orientation \times num\_scale)^{body\_parts}$$

$$= (10^6 \times 10^2 \times 5)^9 = 6 \times 10^{78}!$$

informedia

- Controlled setting:  background subtraction

- Hand initialization, user assisted tracking

- Grouping based body detection:
  - Forsyth, Finding Naked Body, ECCV'96
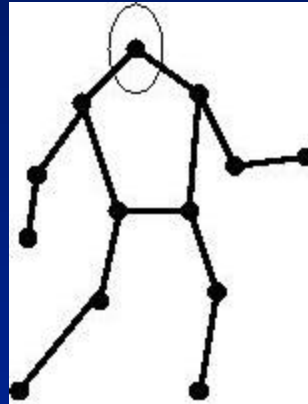  - Felzenszwalb, Huttenlocher, ICCV'99

*informedia*

- Using motion information to find possible locations of joints,

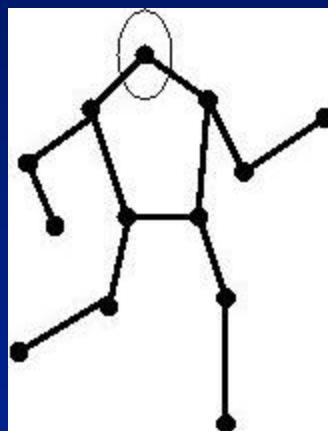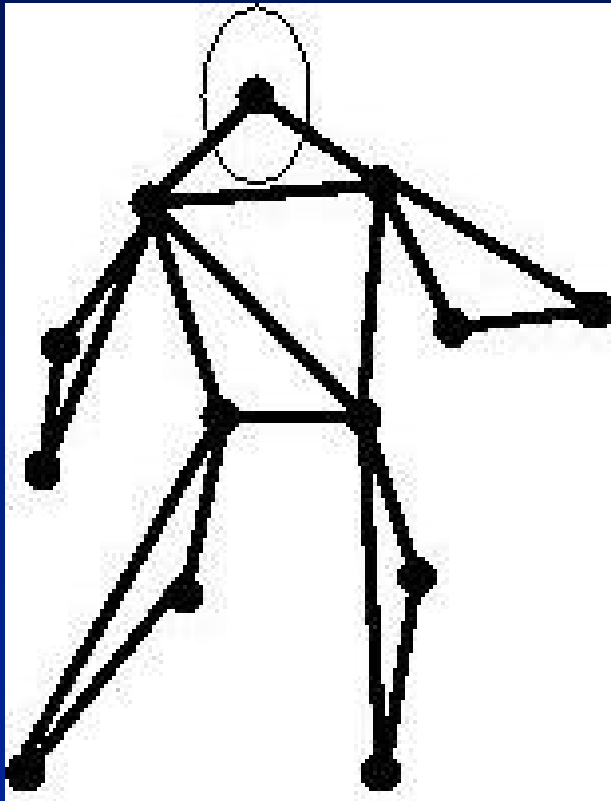- Using MRF inferencing technique for finding the globally optimal body configuration.

*informedia*
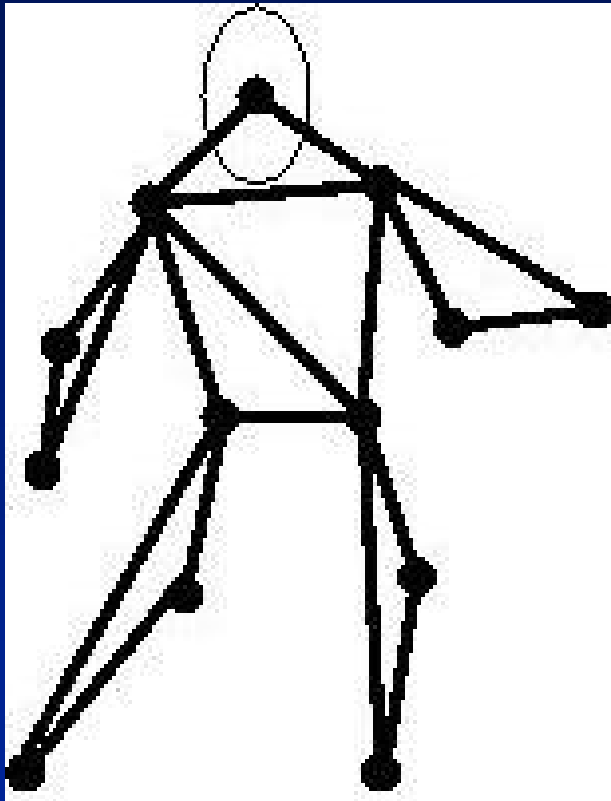
Let $L = (l_1, l_2, ..., l_n)$, be the assignment of each joints

$$E(L) = \sum_{(l_i, l_j, l_k)} E(l_i, l_j, l_k) + \sum_{l_i} E(I, l_i),$$

$$P(L) = \frac{1}{Z} e^{-E(L)}$$

*informedia*

$$E(L) = \sum_{(l_i, l_j, l_k)} E(l_i, l_j, l_k) + \sum_{l_i} E(I, l_i),$$
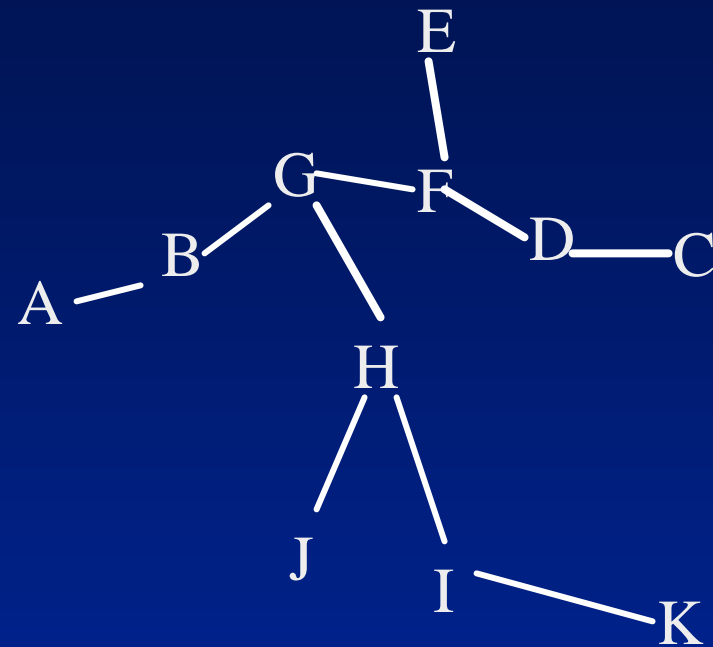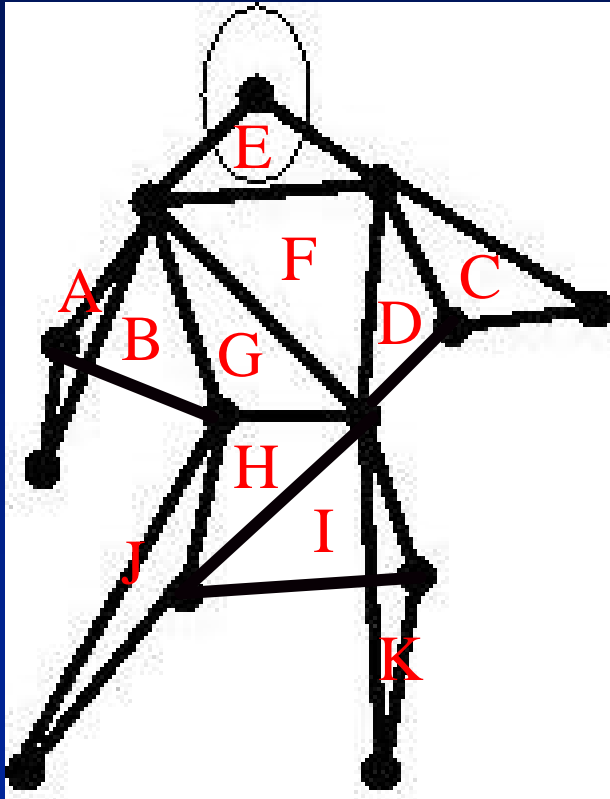
$E(I, l_i) = $ how likely $l_i$ is a joint,

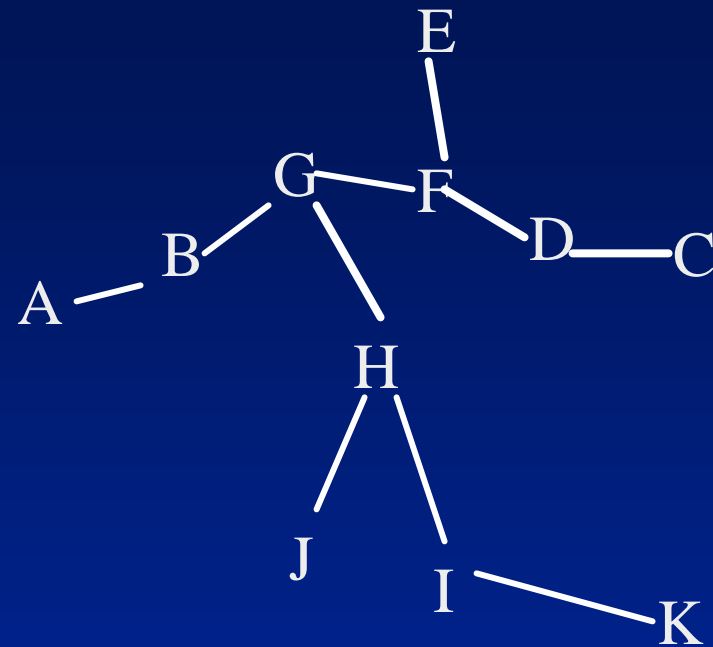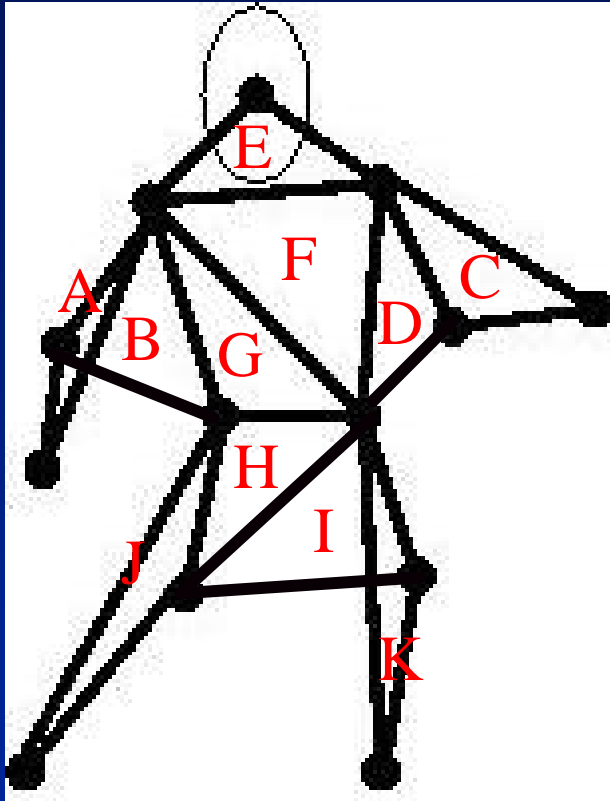$E(l_i, l_j, l_k) = $ how likely is this configurat ion

of three joints.

Depends on :

1) geometrica l relationsh ip,
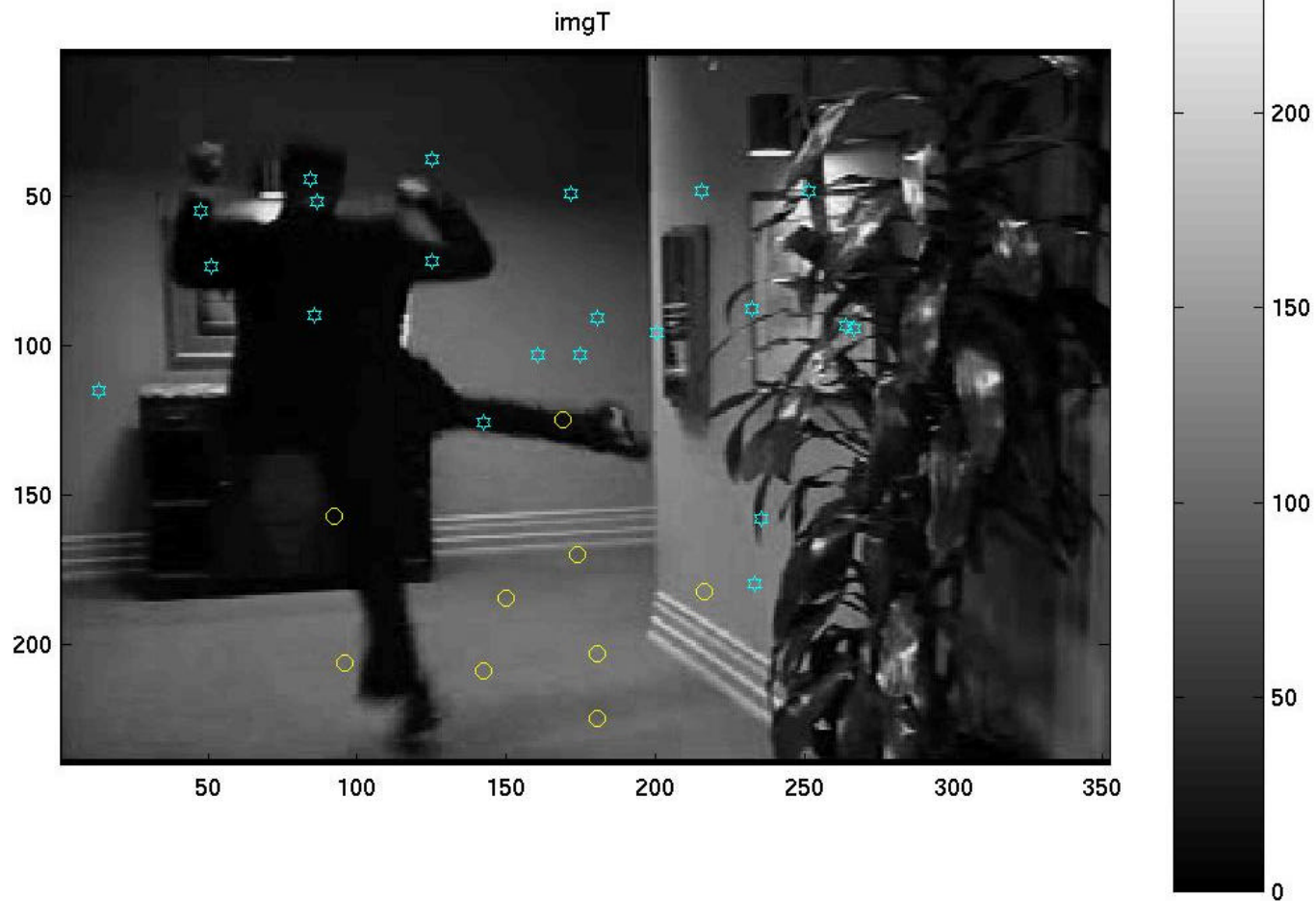
2) relative motion,

3) image informatio n.

Run Dynamic programming on
the clique tree above

Specialize for
indexing
Human and Human
Activities

Intelligent Compact
Video Representation for
Fast trainable query

Data-mining from mapping
Between Video and Text
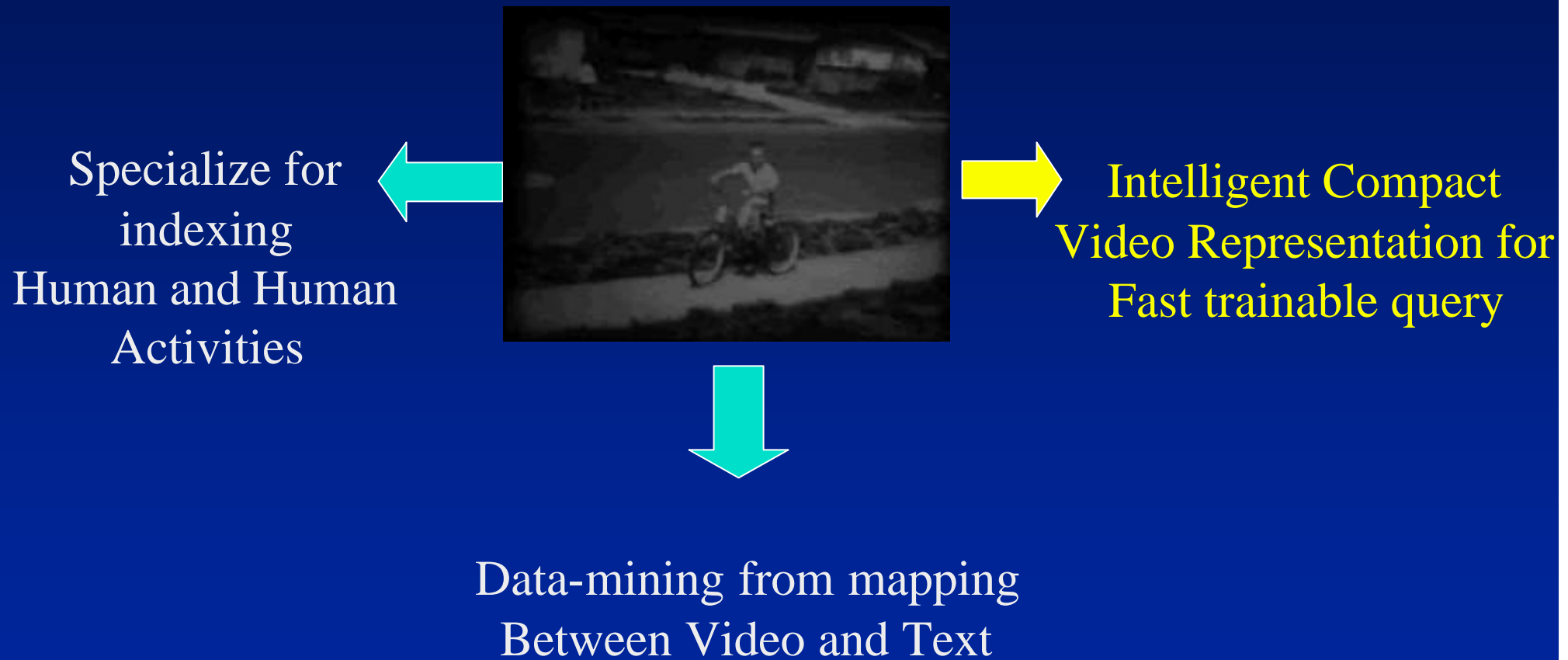
*Original Video Segment*

(3.4Mb)

Compute the transformation
between image I and J, using
affine approximation:

$$J(Ax+D) = I(x)$$
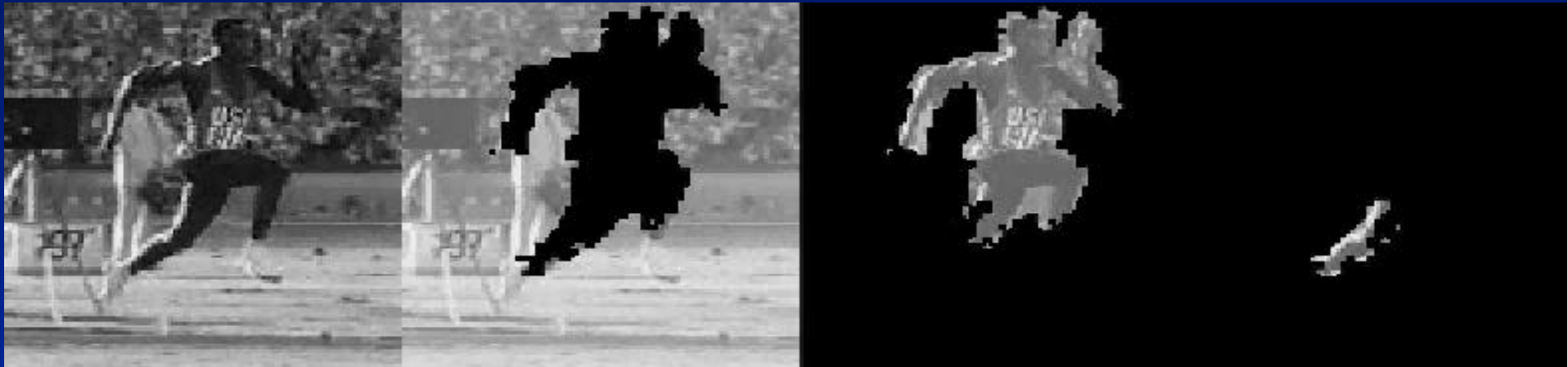
*Original Video Segment*
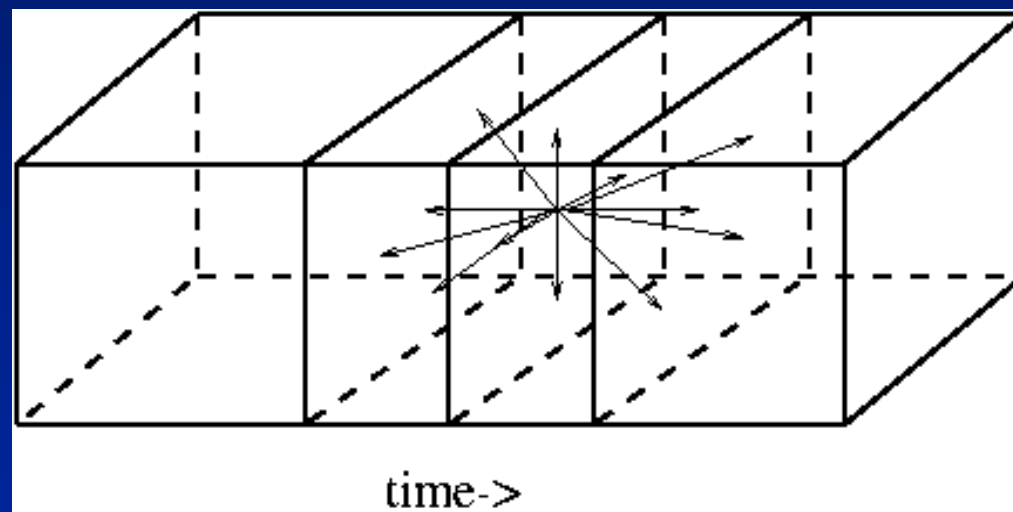
(3.4Mb)

*Panorama Layer*

(35Kb)

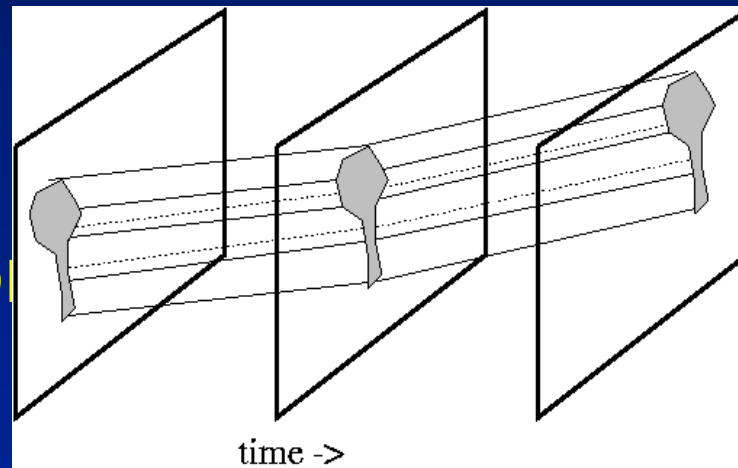- ## Multiple object motion

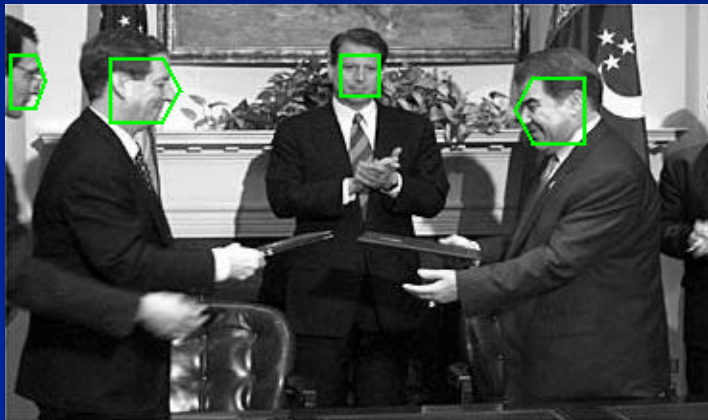- **Networks of spatial-temporal connections:**
  Motion Segmentation with Normalized Cuts



time->
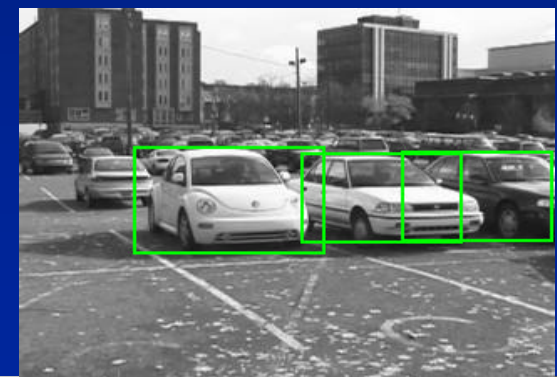
- ## Motion "proto-volume" in space-time



time ->

- ## Group co

- **Representation objects/scene at multiple level of abstraction**
  - Low level motion, texture description
  - Mid level object segmentation
  - High level object type and instances

*informedia*

# Learning Seg. With Random Walk

**Specialize for indexing Human and Human Activities**

**Intelligent Compact Video Representation for Fast trainable query**

**Data-mining from mapping Between Video and Text**

- Looking for cell phone images

- **Keyword annotation is not sufficient**
  - Too many objects
  - Disagreement of keywords

- **Move towards more powerful image-query system in Video**



Specialize for
indexing
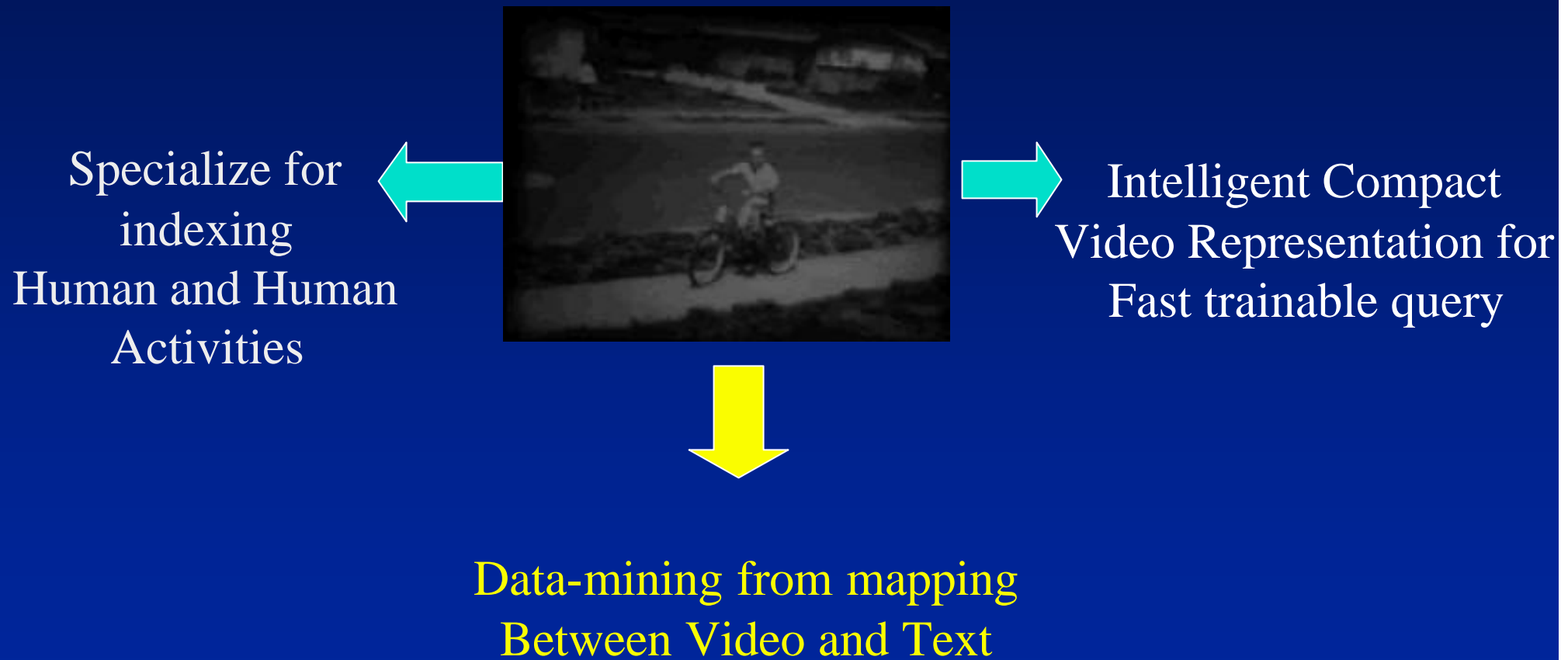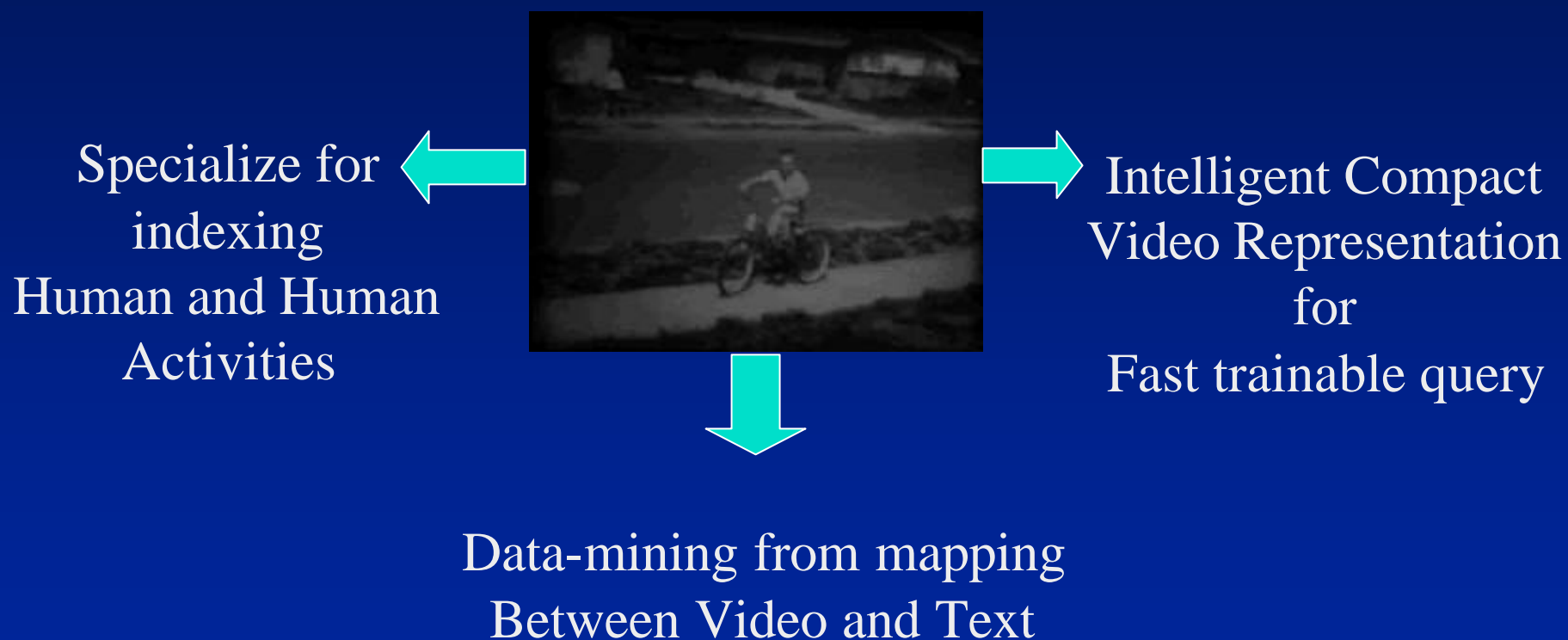Human and Human
Activities

Intelligent Compact
Video Representation
for
Fast trainable query

Data-mining from mapping
Between Video and Text

*informedia*

- **Move towards more powerful image-query system in Video**



Specialize for
indexing
Human and Human
Activities

Intelligent Compact
Video Representation
for
Fast trainable query

Data-mining from Co-relation
Between Video and Text